

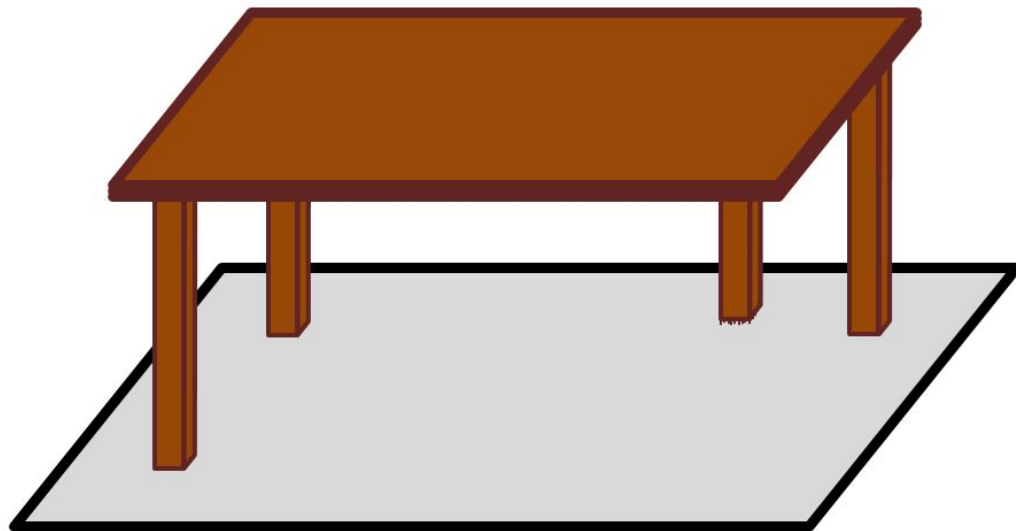
A new paradigm for computer vision based on compositional representation

Vinjai Vale
Mentor: Kevin Ellis

Computer vision today

- Image processing, analysis, understanding
- State of the art: deep convolutional neural network (CNN)
- CNNs excel at *classification*, struggle with *representation*

Example



Object compositionality

Represent objects recursively through their components and relations

Enables complex human-level visual reasoning

CNN drawbacks

CNNs learn fuzzy patterns on textures — poor spatial understanding

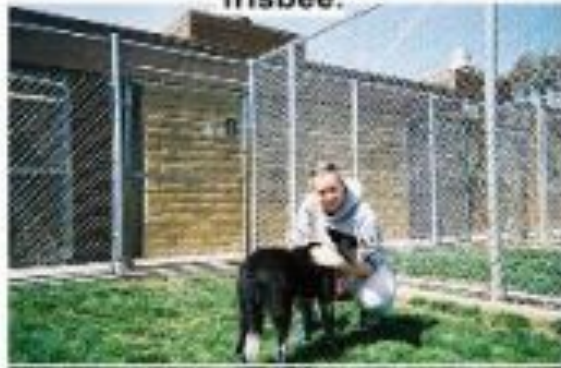


<https://rocknrollnerd.github.io/ml/2015/05/27/leopard-sofa.html>
<https://sleezybarbhorsetwear.com/customers-page/>



T. Salimans, I. J. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training GANs. CoRR, abs/1606.03498, 2016.

A dog is jumping to catch a frisbee.



A refrigerator filled with lots of food and drinks.



A yellow school bus parked in a parking lot.



Research goal

- Engineer domain and dataset to isolate the problem of compositionality
- Develop algorithms, techniques to solve representation tasks on dataset

Primitive elements and ShapeWorld

A small number of *primitive elements* (PEs) compose all objects.

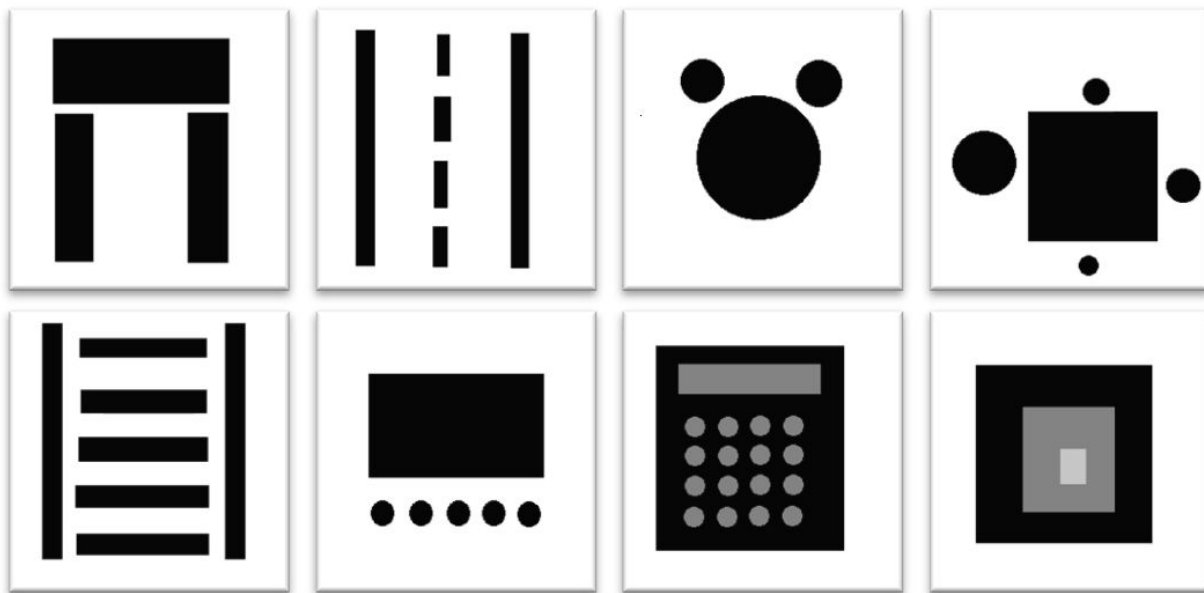
Biederman 1987: 36 *geons* are the PEs of human vision in 3D world

Key intuition:

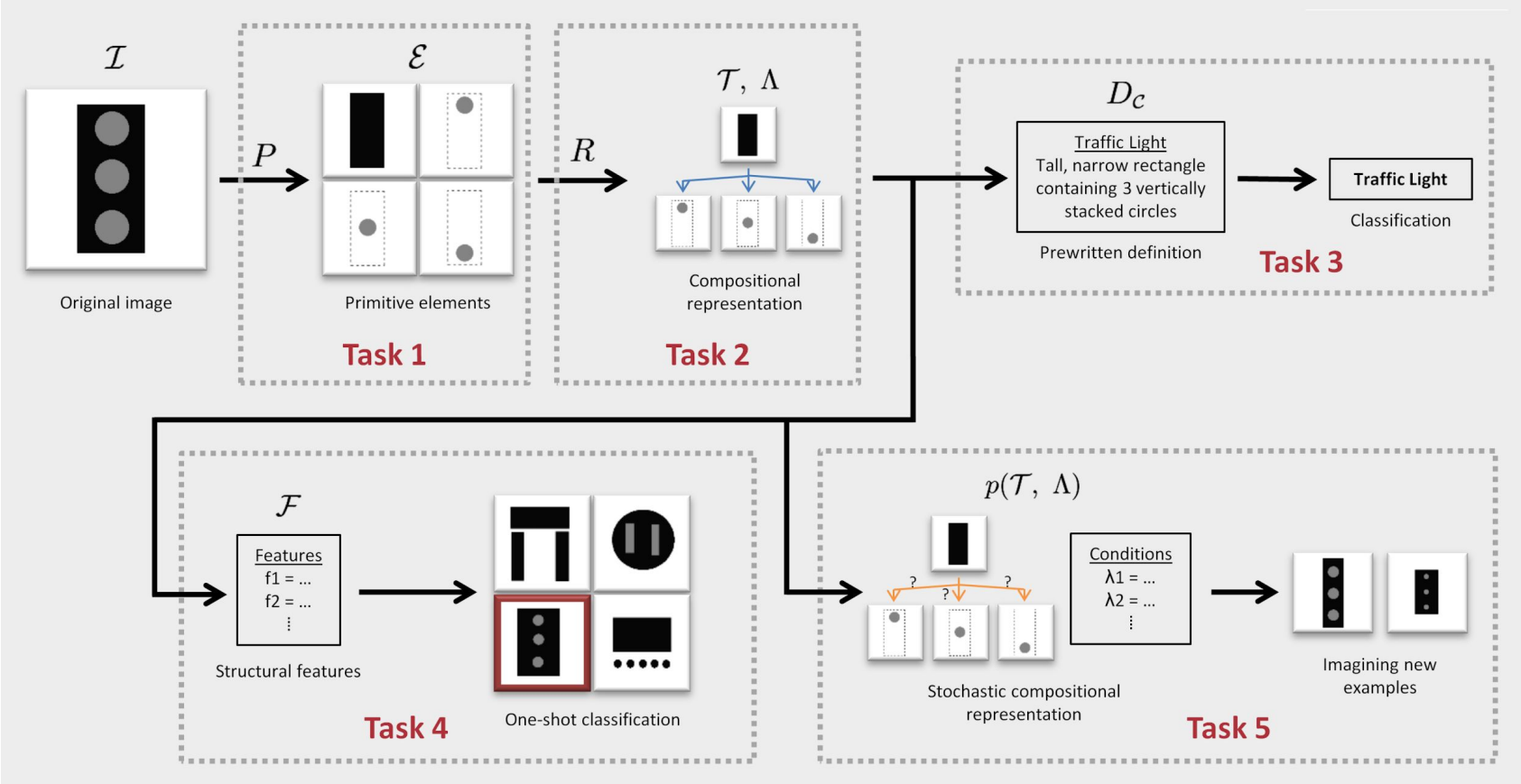
- **3D world with 36 PEs is hard...**
- **Instead, solve vision in 2D world with only 2 PEs!**

Primitive elements and ShapeWorld

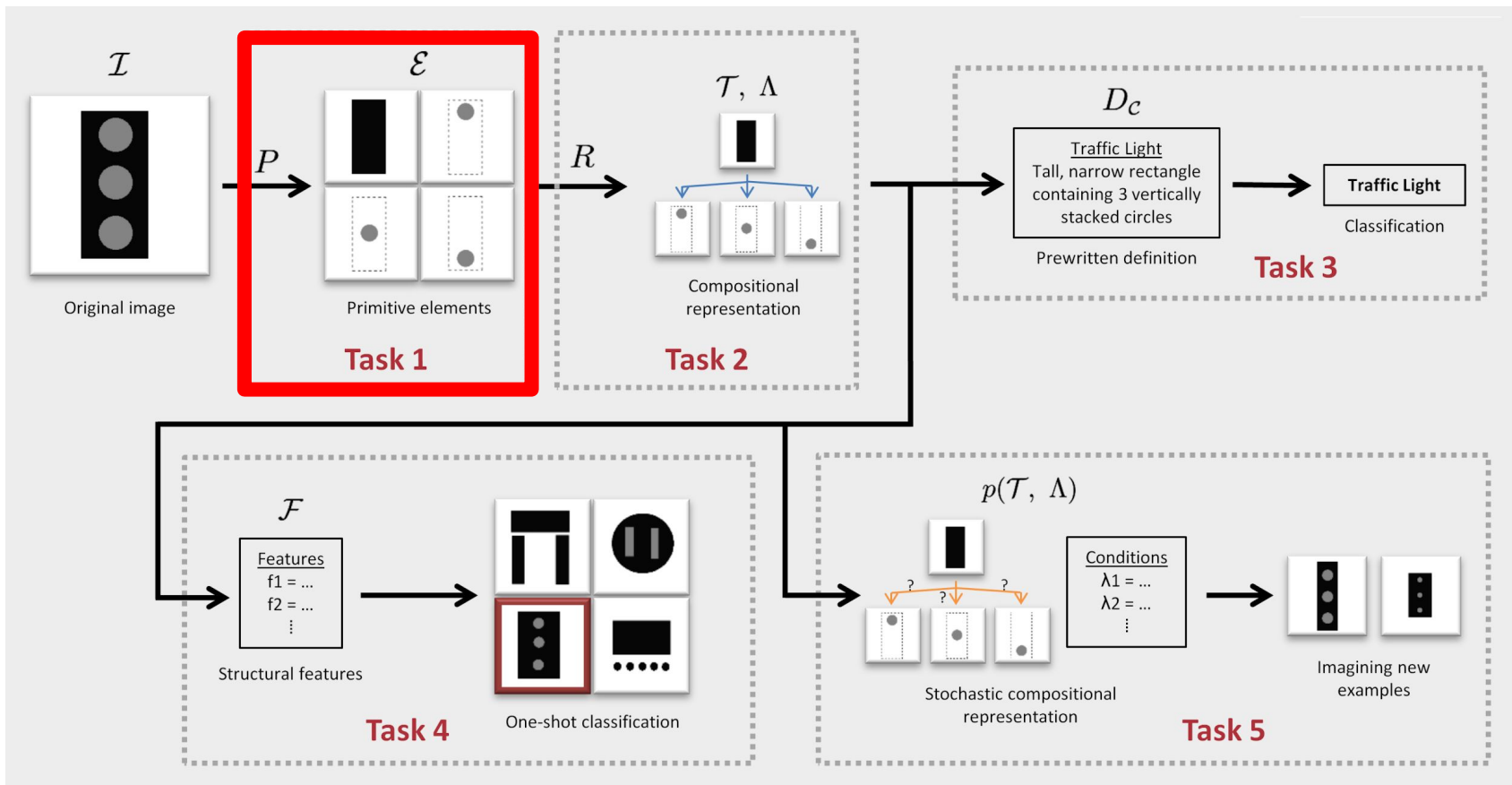
ShapeWorld: 2D dataset composed of circles and rectangles



Five Tasks



Task 1



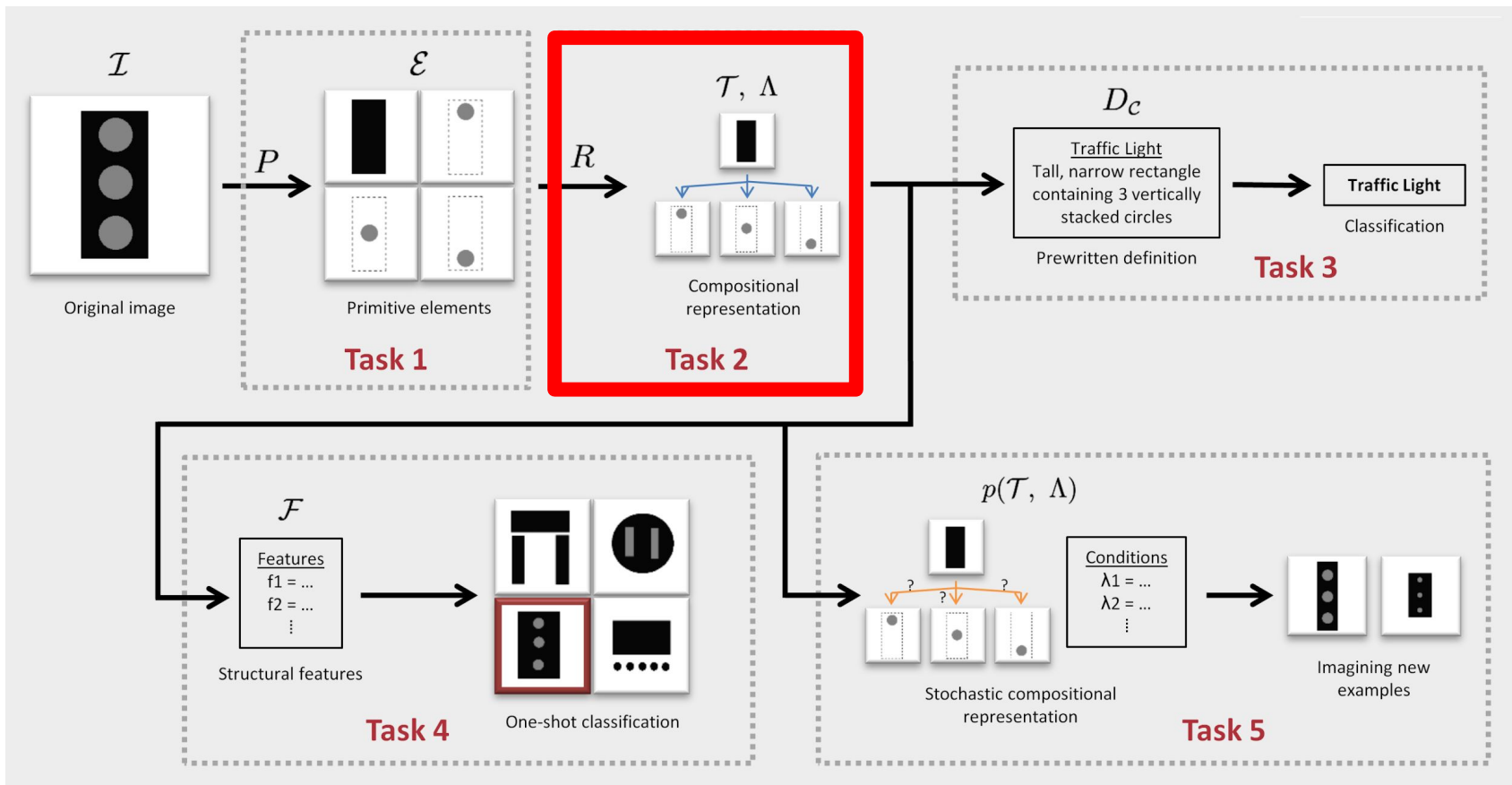
Task 1: PE decomposition

Represent an image in terms of a *graphics program*

“Draw circle here with these coordinates,” “draw rectangle there with those coordinates,” etc.

Image processing methods - OpenCV

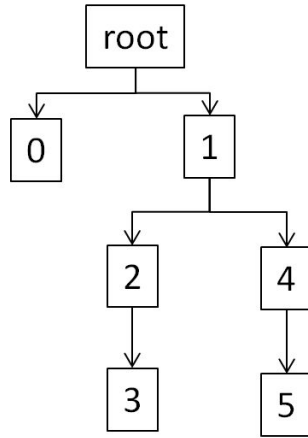
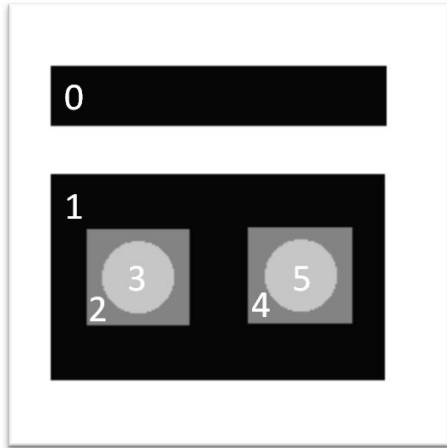
Task 2



Task 2: Compositional representation

Two-part data structure to encompass compositional relationships

Augmented Primitive Element Tree (APET)

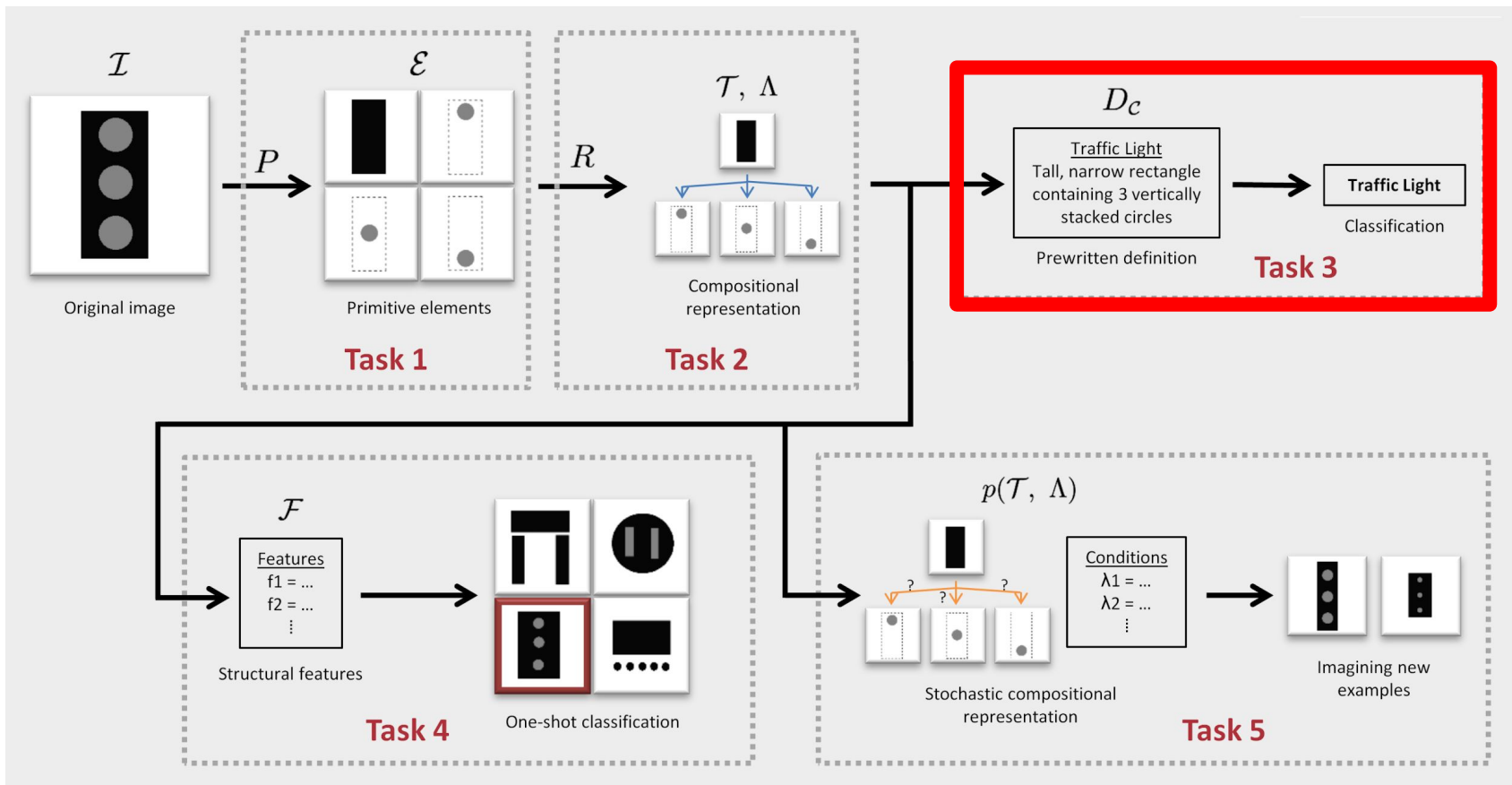


Coincidence List

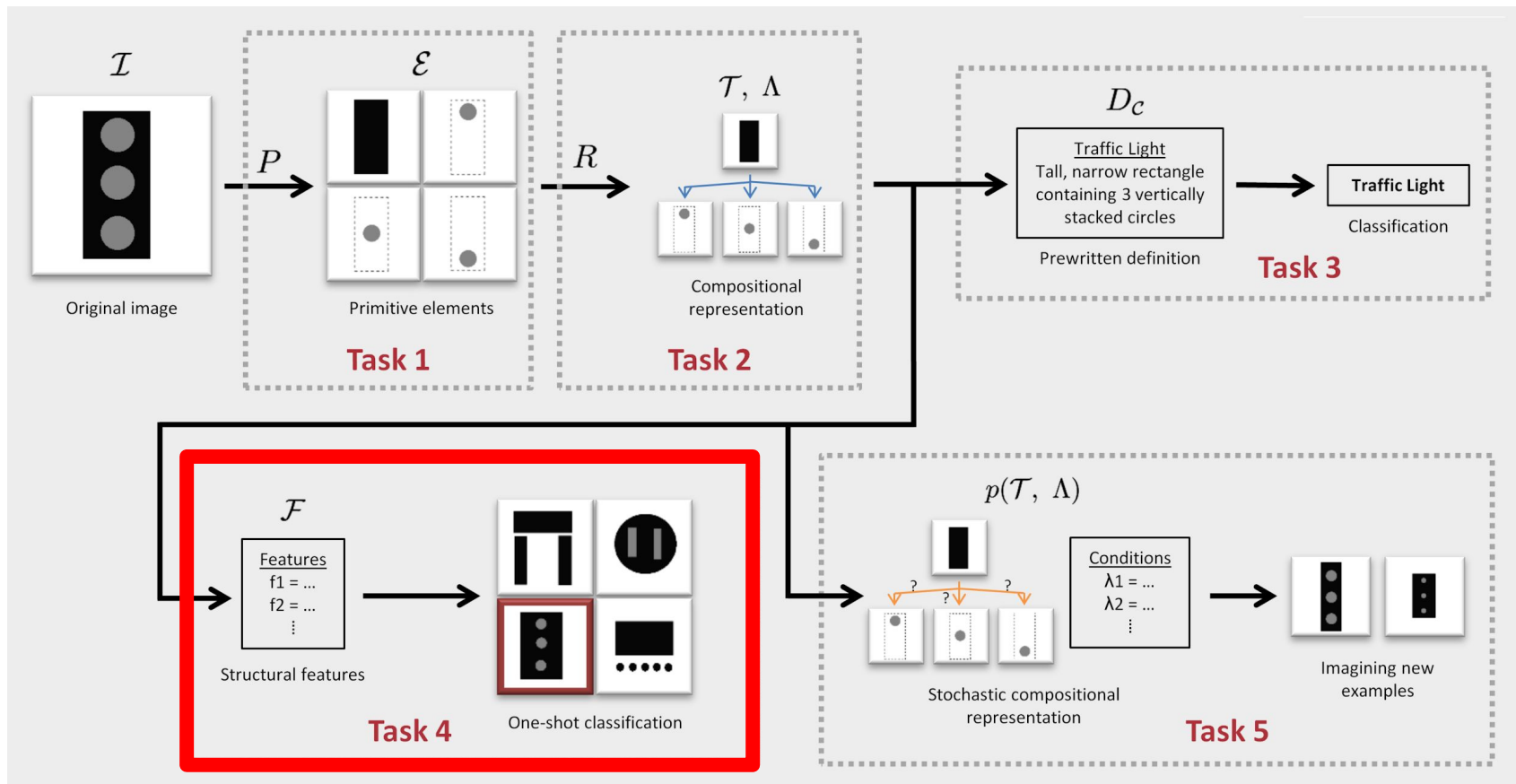
Coincidences likely part of underlying concept

- Rows
- Clusters
- Grids
- Radial arrangement
- Etc.

Task 3



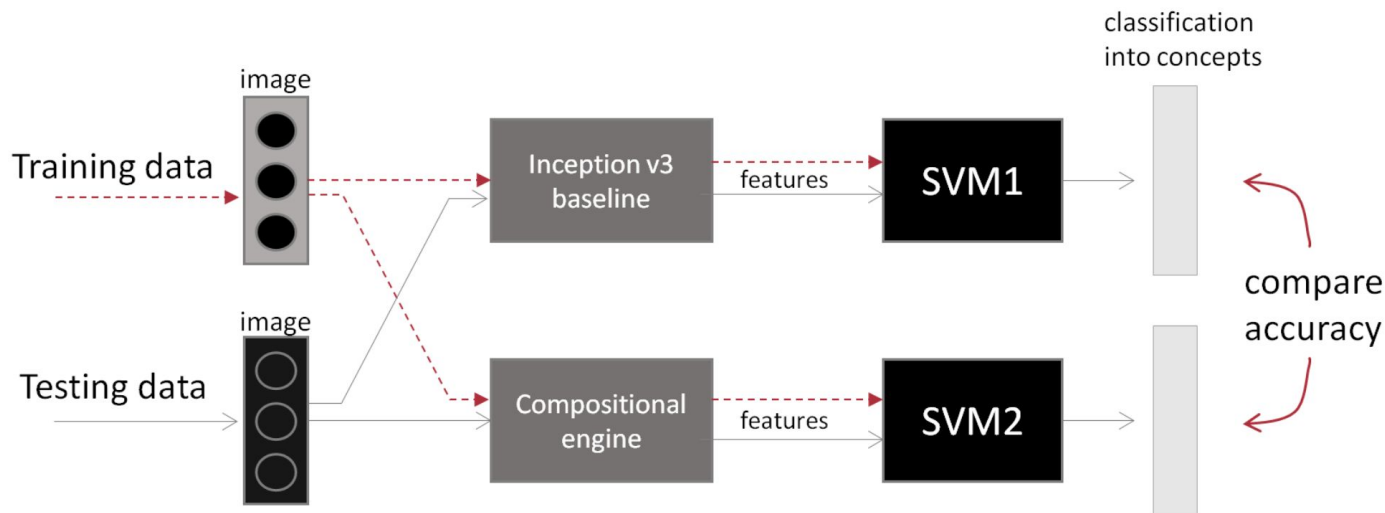
Task 4



Task 4: Baseline comparison

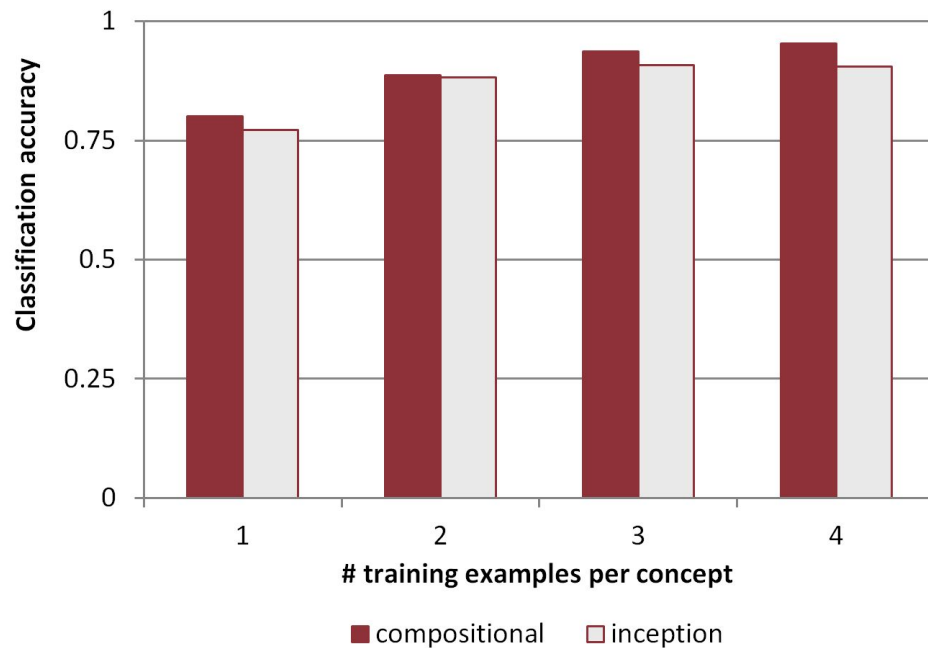
Extracted 34 features from Task 2 compositional representation, compared against 2048 features from Google Inception-v3 CNN

The task: learn to classify a ShapeWorld concept only given 1-4 examples



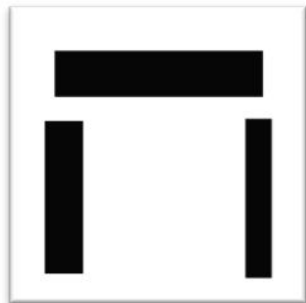
Task 4 results

Compositional vs. Google Inception-v3 Features

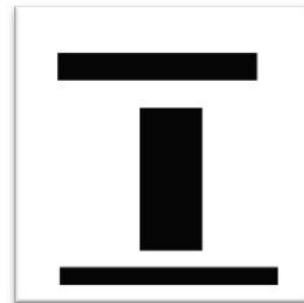
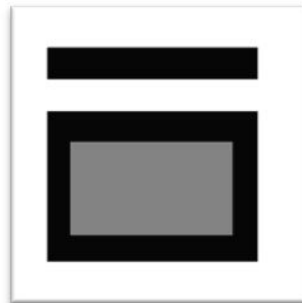


Task 4 closer inspection

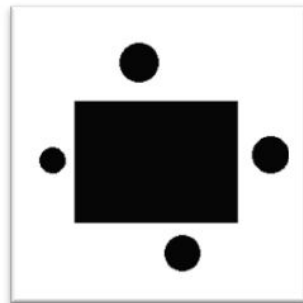
Inception-v3



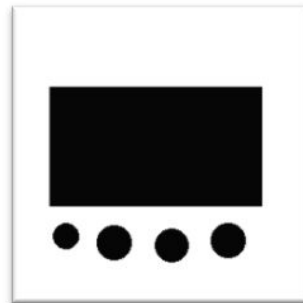
mistaken for



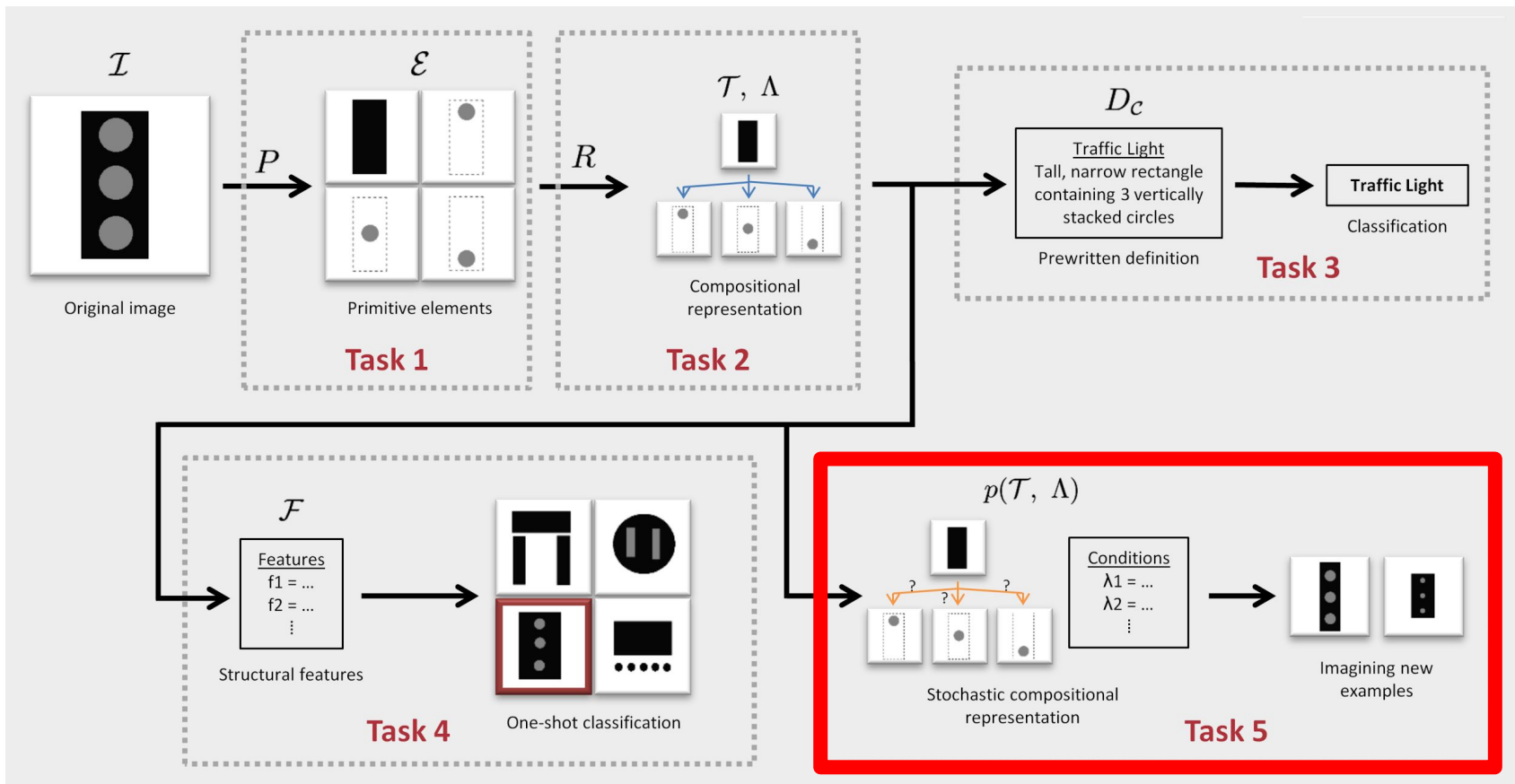
Compositional features



mistaken for



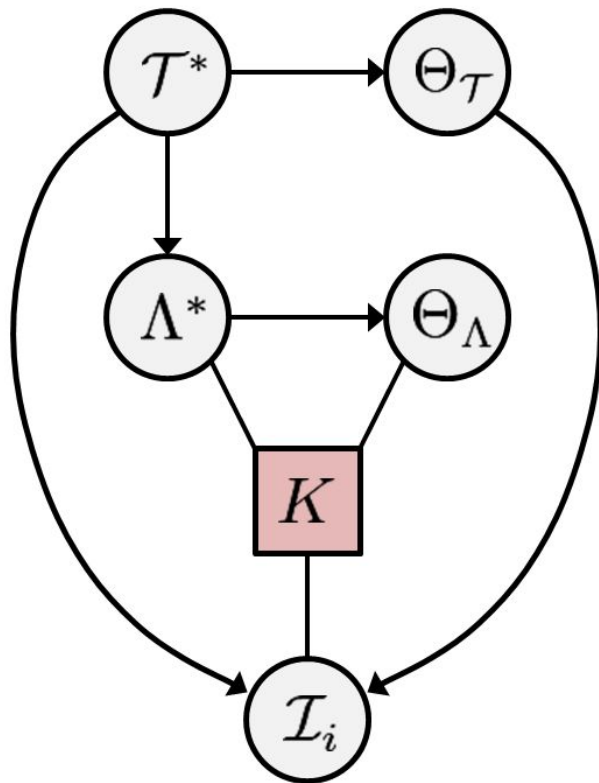
Task 5



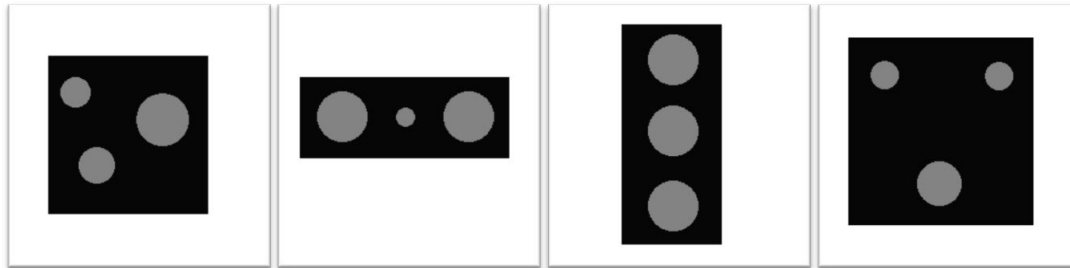
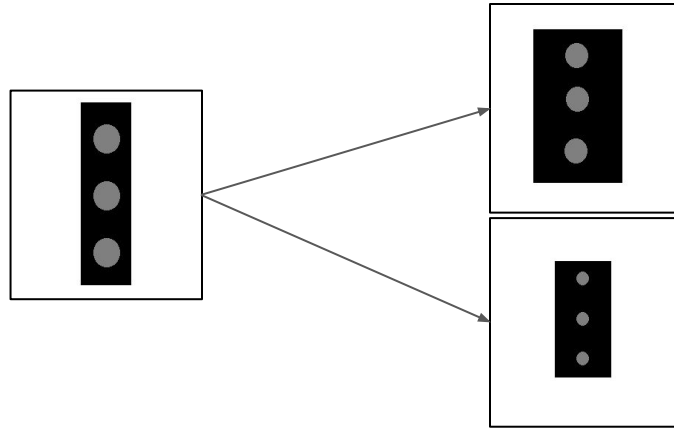
Task 5: Learning a generative model

Mathematical model for inferring a probability distribution $P(I, S_c)$: probability that image I belongs to concept C

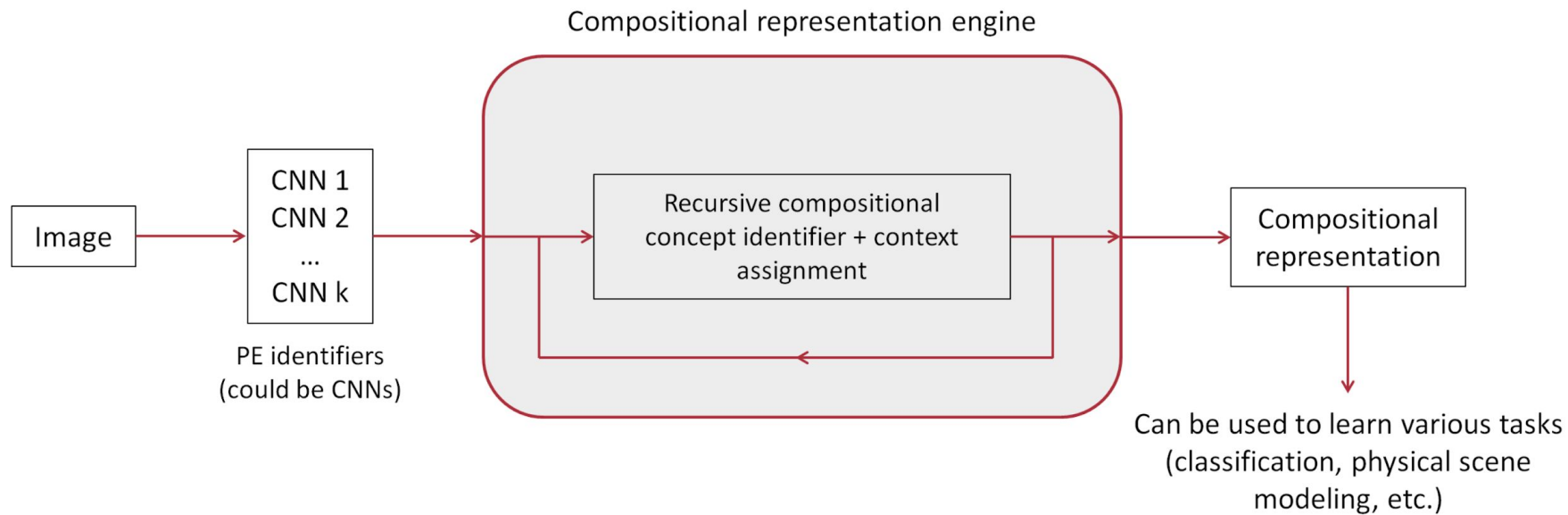
Task 5: Learning a generative model



Task 5: Learning a generative model



Beyond



Conclusion

1. ShapeWorld dataset following Recognition-by-Components paradigm
2. Generate compositional representations of images (APET, Coincidence List)
3. Compositional representation feature set outperforms state-of-the-art CNN in object compositionality tasks
4. Mathematical, generalizable probabilistic approach to learning stochastic compositional representations

Goal: AI vision systems that are faster, safer, and closer to human vision.

Acknowledgements

Mentor Kevin Ellis @ MIT Brain and Cognitive Sciences

MIT PRIMES Program

Szczesny Kaminski and Sean Campbell, my teachers, for facilitating independent studies in machine learning